

# Measuring Bayesian (and variational Bayesian) robustness

---

Wow wow

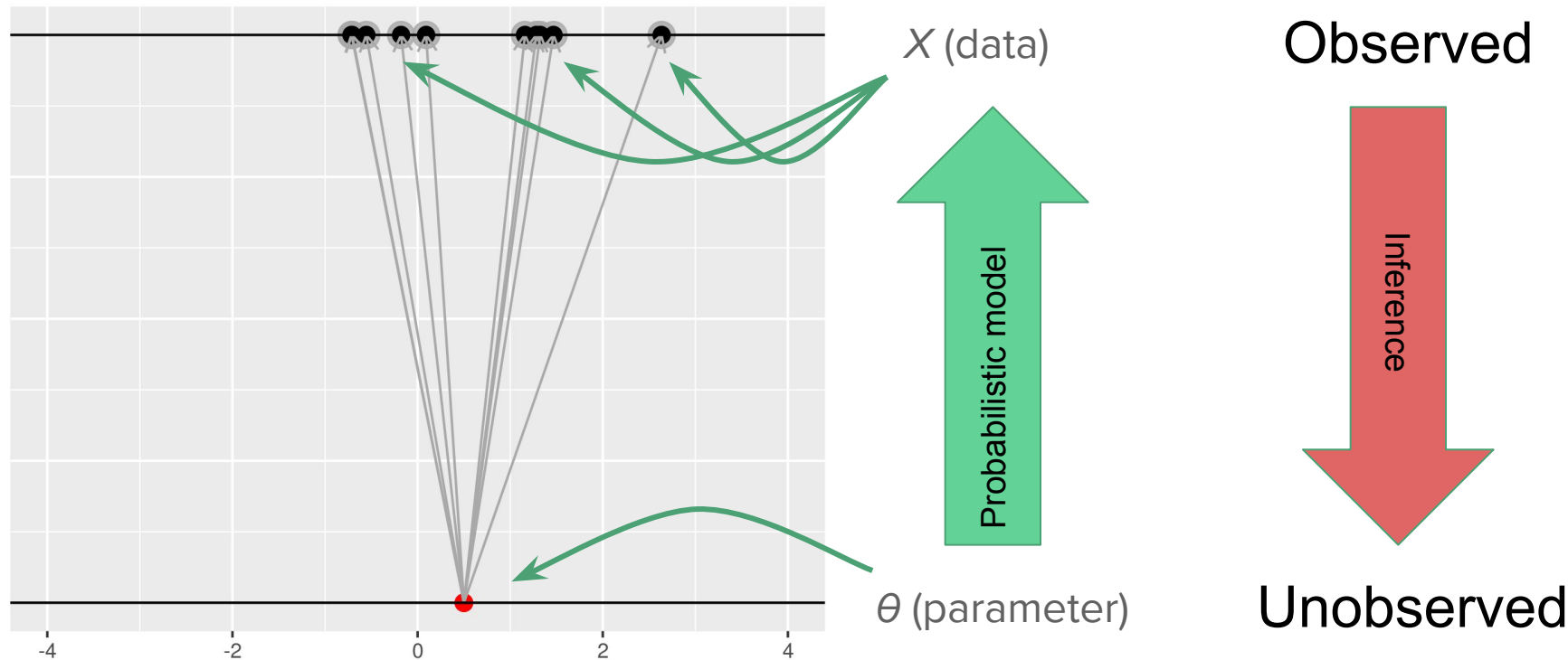
# Outline

- What is Bayesian statistics?
  - ...as opposed to frequentist statistics?
  - Bayes rule and posterior approximation
  - Linear regression example
- What do we mean by Bayesian robustness?
  - Local robustness
  - Motivating example: Bayesian “leverage scores”
  - Sensitivity is covariance is sensitivity
  - MCMC for approximate leverage scores?
- What is variational Bayes?
  - The closest distribution in KL divergence
  - Mean field approximations
- How do you calculate robustness to function-valued perturbations?
  - The influence function
  - Worst-case perturbations

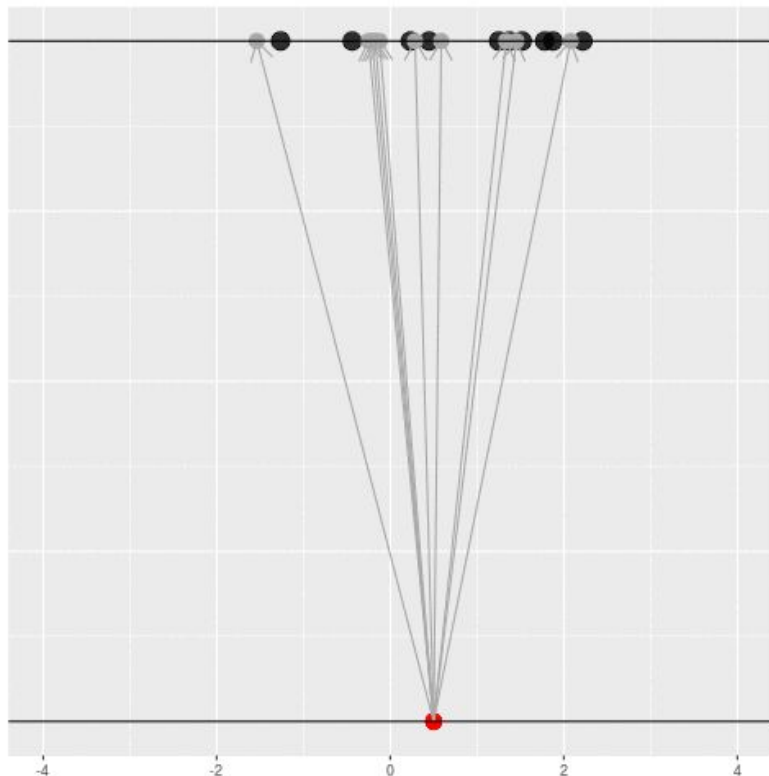
# What is Bayesian statistics?

---

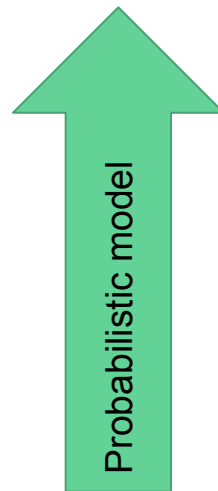
# Parameters and data



# Frequentist approach



$X$  (data)



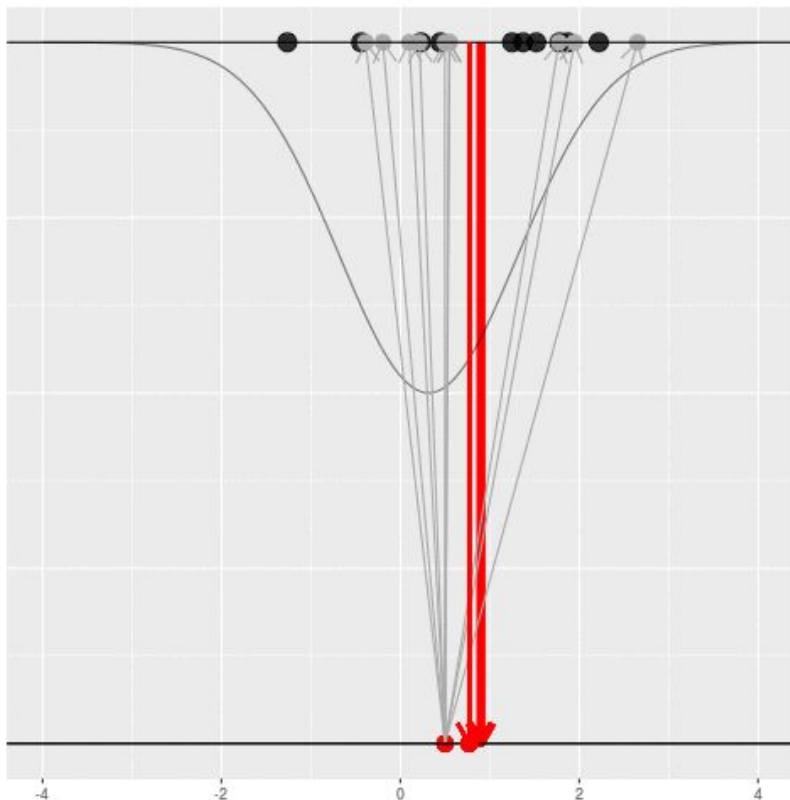
$\theta$  (parameter)

## Frequentist idea:

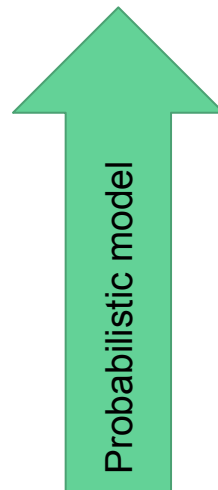
We got the parameter indicated by the red dot and saw the dataset in black.

But the same parameter could have given us lots of other datasets.

# Frequentist approach



$X$  (data)



$\theta$  (parameter)

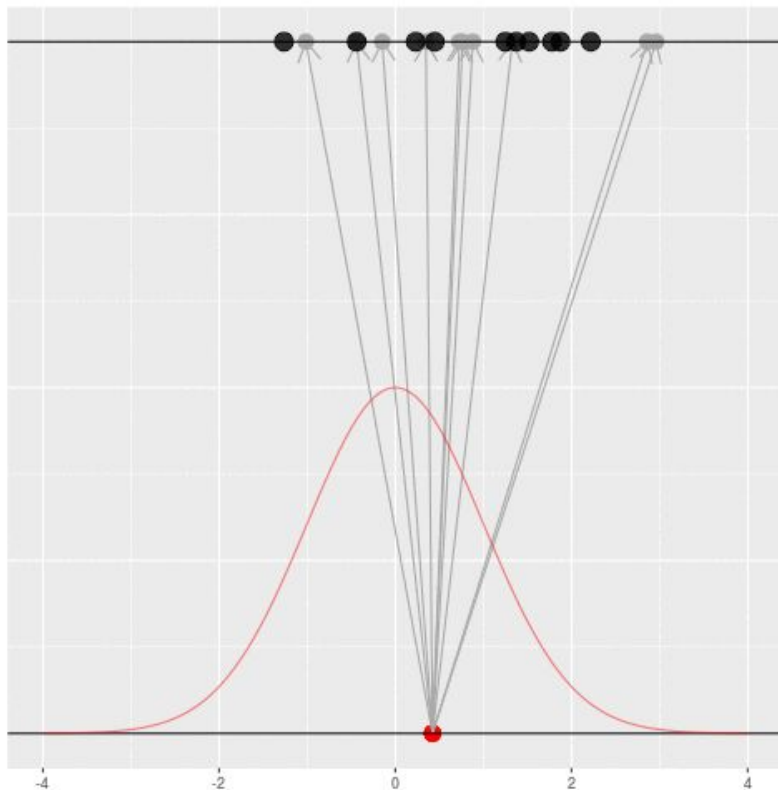
## Frequentist idea:

For each dataset, we might pick some summary function and call it an “estimate”.

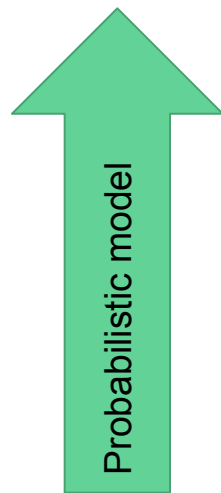
It will be different each time because the data will be different each time.

We hope the estimate is usually near the true parameter in some sense.

# Bayesian approach



$X$  (data)



$\theta$  (parameter)

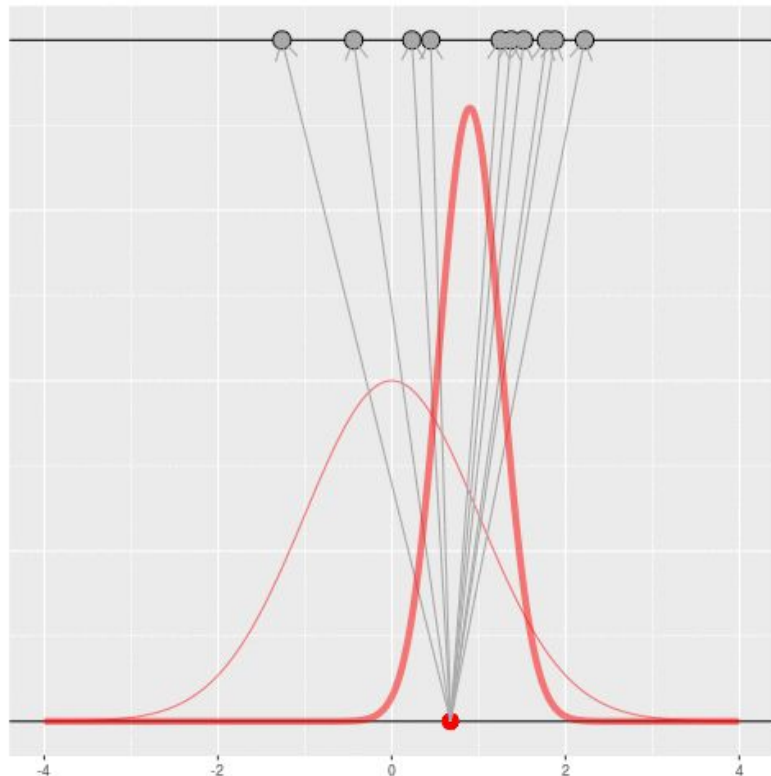
## Bayesian idea:

Let's imagine that the universe generates datasets by

- First drawing a parameter from some prior distribution
- and then drawing a dataset from that parameter.

Sometimes we would get the dataset we saw. Mostly we wouldn't.

# Bayesian approach



$X$  (data)

Probabilistic model

$\theta$  (parameter)

## Bayesian idea:

Suppose we draw a bunch of parameters and datasets, and then throw out every pair where the data doesn't match what we observed.

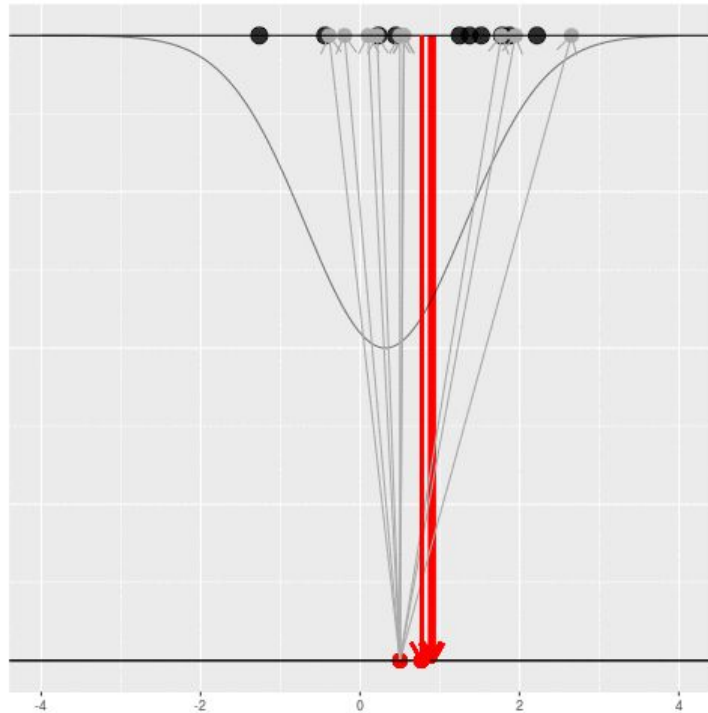
The distribution of the parameters that are left represents which parameters could have given us the dataset we saw.

We hope the prior is reasonable and the model is accurate.

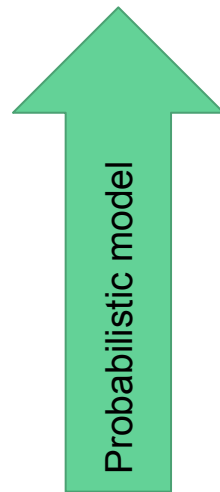


Frequentist: the parameter is fixed, data is random.

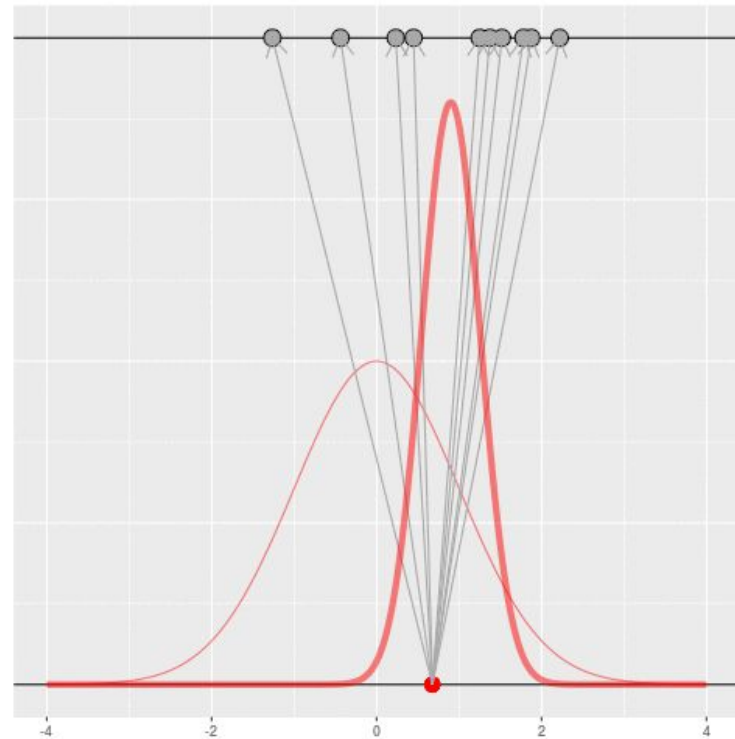
Bayes: the data is fixed, the parameter is random.



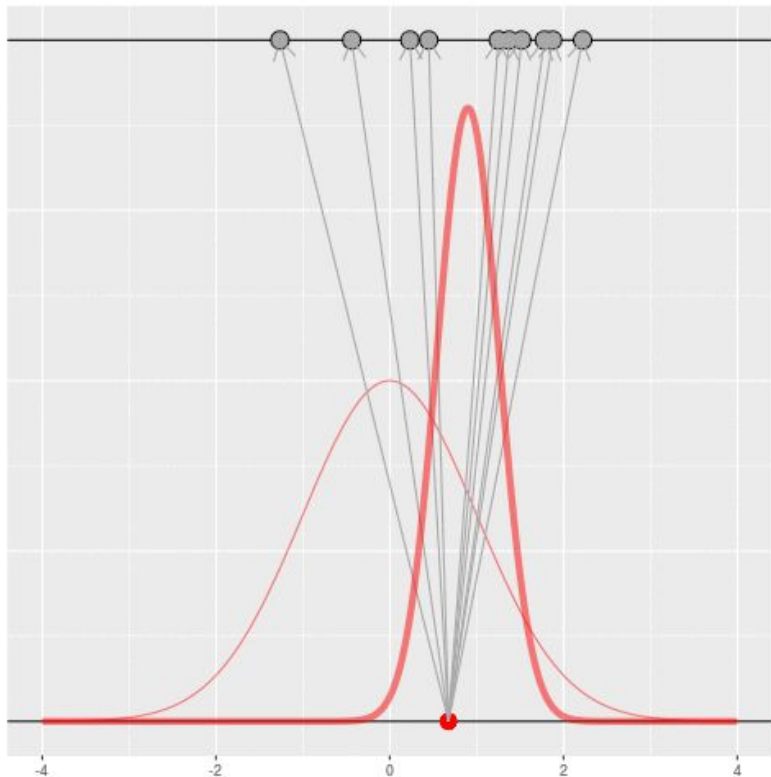
$X$  (data)



$\theta$  (parameter)



# Bayesian approach



Of course, in practice you don't usually generate parameters and data hoping to get your original dataset.

Instead, you use Bayes' rule:

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

This is intractable in general (the denominator is a problem). Turn to approximations schemes like MCMC, variational Bayes, &c.

# Linear regression example

$$p(y_i|\beta) = \mathcal{N}(x_i^T \beta, \sigma^2)$$

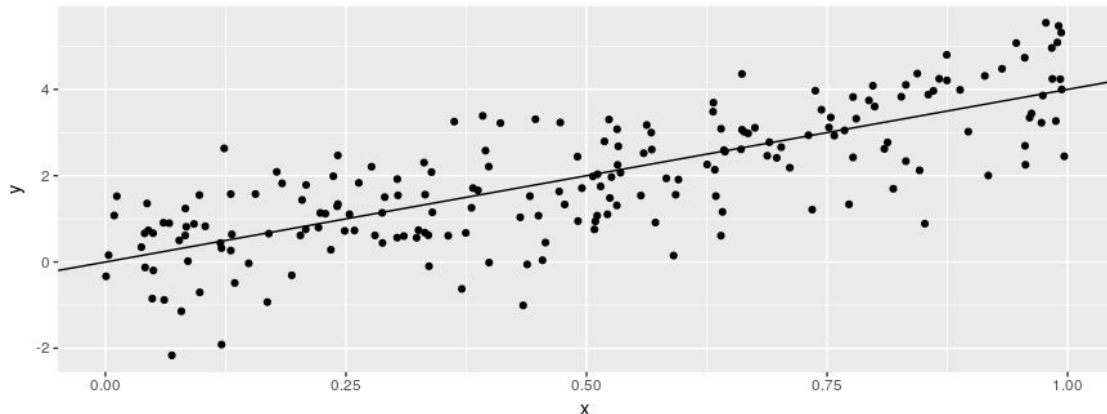
$y_i$  = Observed scalar (data)

$\beta$  = Unobserved  $p$  – dimensional column vector (parameter)

$x_i$  = Observed  $p$  – dimensional column vector

$\sigma^2$  = Observed (known) variance

} Could think of these as “hyperparameters” -- they are not modeled as random, we assume we just know them



1d example:

We observe the points and want to infer the line.

# Linear regression example

We observe  $n$  data points, so in matrix notation:

$$p(Y|\beta) = \mathcal{N}(X\beta, \sigma^2 I_n)$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \text{n-dimensional column vector}$$

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = \text{n by p matrix}$$

# Linear regression example

We need to assume a prior.

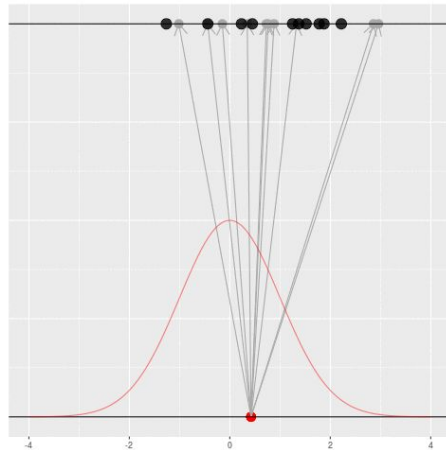
A convenient (“conjugate”) form is a normal prior:

$$p(Y|\beta) = \mathcal{N}(X\beta, \sigma^2 I_n)$$

$$p(\beta) = \mathcal{N}(0, \sigma_0^2 I_p)$$

Then:

$$\begin{aligned}\log p(Y, \beta) &= \log(p(Y|\beta) p(\beta)) \\ &= -\frac{1}{2}\sigma^{-2} (Y - X\beta)^T (Y - X\beta) - \frac{1}{2}\sigma_0^{-2} \beta^T \beta + C\end{aligned}$$



# Linear regression example

$$\begin{aligned}\log p(Y, \beta) &= \log(p(Y|\beta) p(\beta)) \\&= -\frac{1}{2}\sigma^{-2} (Y - X\beta)^T (Y - X\beta) - \frac{1}{2}\sigma_0^{-2} \beta^T \beta + C \\&= -\frac{1}{2}\sigma^{-2} \beta^T X^T X \beta + \sigma^{-2} Y^T X \beta - \frac{1}{2}\sigma_0^{-2} \beta^T \beta + C \\&= -\frac{1}{2}\text{trace}((\sigma^{-2} X^T X + \sigma_0^{-2} I_p) \beta \beta^T) + \sigma^{-2} Y^T X \beta + C\end{aligned}$$

Recall that, for a generic multivariate normal distribution,

$$\begin{aligned}\log \mathcal{N}(\beta; \mu, \Sigma) &= -\frac{1}{2} (\beta - \mu)^T \Sigma^{-1} (\beta - \mu) + C \\&= -\frac{1}{2}\text{trace}(\Sigma^{-1} \beta \beta^T) + \mu^T \Sigma^{-1} \beta + C\end{aligned}$$

# Linear regression posterior

$$\begin{aligned}\log p(Y, \beta) &= \log(p(Y|\beta) p(\beta)) \\ &= -\frac{1}{2} \text{trace}((\sigma^{-2} X^T X + \sigma_0^{-2} I_p) \beta \beta^T) + \sigma^{-2} Y^T X \beta + C\end{aligned}$$

Recall that, for a generic multivariate normal distribution,

$$\begin{aligned}\log \mathcal{N}(\beta; \mu, \Sigma) &= -\frac{1}{2} (\beta - \mu)^T \Sigma^{-1} (\beta - \mu) + C \\ &= -\frac{1}{2} \text{trace}(\Sigma^{-1} \beta \beta^T) + \mu^T \Sigma^{-1} \beta + C\end{aligned}$$

Therefore the posterior is given by:

$$p(\beta|Y) = \mathcal{N}\left((\sigma^{-2} X^T X + \sigma_0^{-2} I_p)^{-1} \sigma^{-2} X^T Y, (\sigma^{-2} X^T X + \sigma_0^{-2} I_p)^{-1}\right)$$

# Linear regression posterior

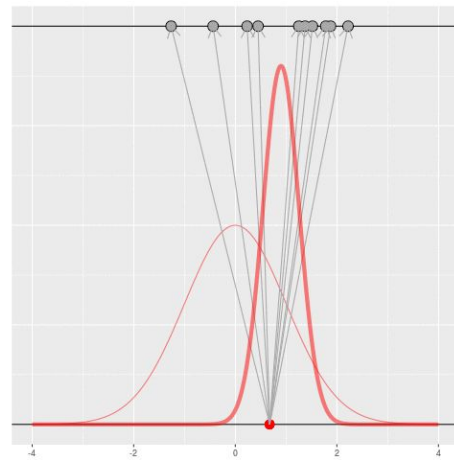
$$p(\beta|Y) = \mathcal{N}\left((\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1} \sigma^{-2}X^T Y, (\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1}\right)$$

$$\mathbb{E}[\beta|Y] = (\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1} \sigma^{-2}X^T Y$$

$$\text{Cov}(\beta|Y) = (\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1}$$

Taking  $\sigma_0 \rightarrow \infty$  gives the OLS solution:

$$\begin{aligned}\mathbb{E}[\beta|Y] &= (\sigma^{-2}X^T X)^{-1} \sigma^{-2}X^T Y \\ &= (X^T X)^{-1} X^T Y\end{aligned}$$





# Bayesian robustness

---

# Bayesian robustness

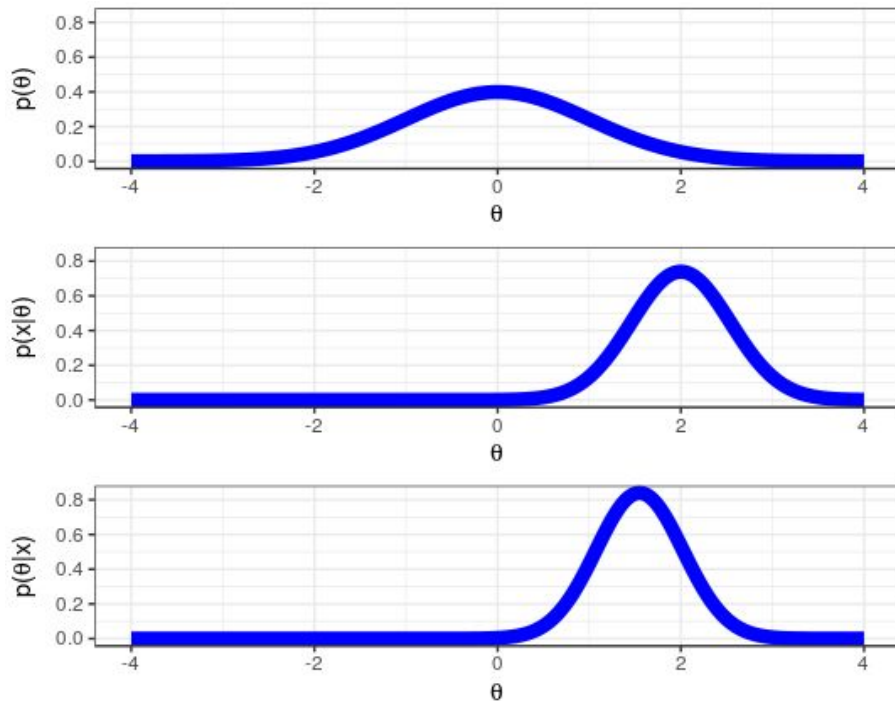
$$p(\theta|\alpha) = \text{Prior}$$

&

$$p(x|\theta) = \text{Likelihood}$$

BAYES

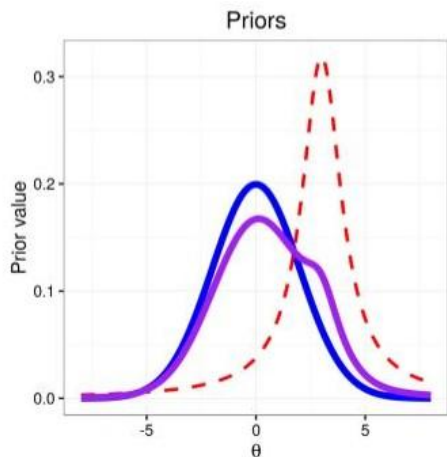
$$p(\theta|x, \alpha) = p_{\alpha}^x(\theta) = \text{Posterior}$$



## Bayesian Data Analysis

Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin

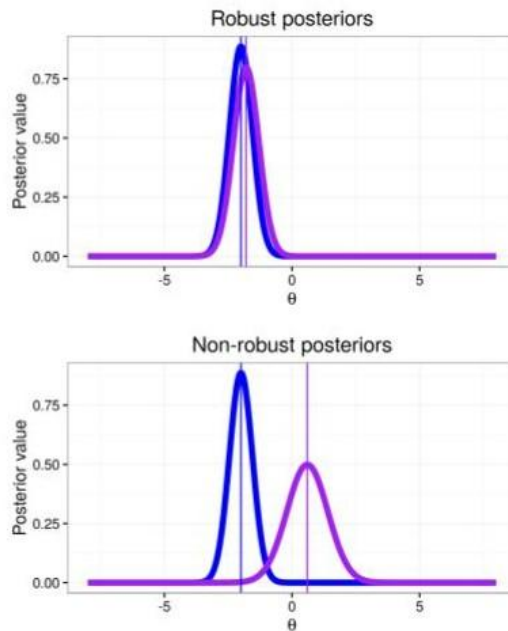
# How do we measure robustness?



Prior type

- $p_0$
- $p_c$
- $p(\theta|\epsilon)$

Data and  
Bayes' Rule



Local sensitivity:

Posterior expectation of interest

$$\frac{d\mathbb{E}[p_{\alpha}^x(\theta)]}{d\alpha} \quad \text{or} \quad \frac{d\mathbb{E}[p_{\alpha}^x(\theta)]}{dx}$$

Perturbation size (of  
hyperparameters or data)

**Local Robustness in Bayesian Analysis**

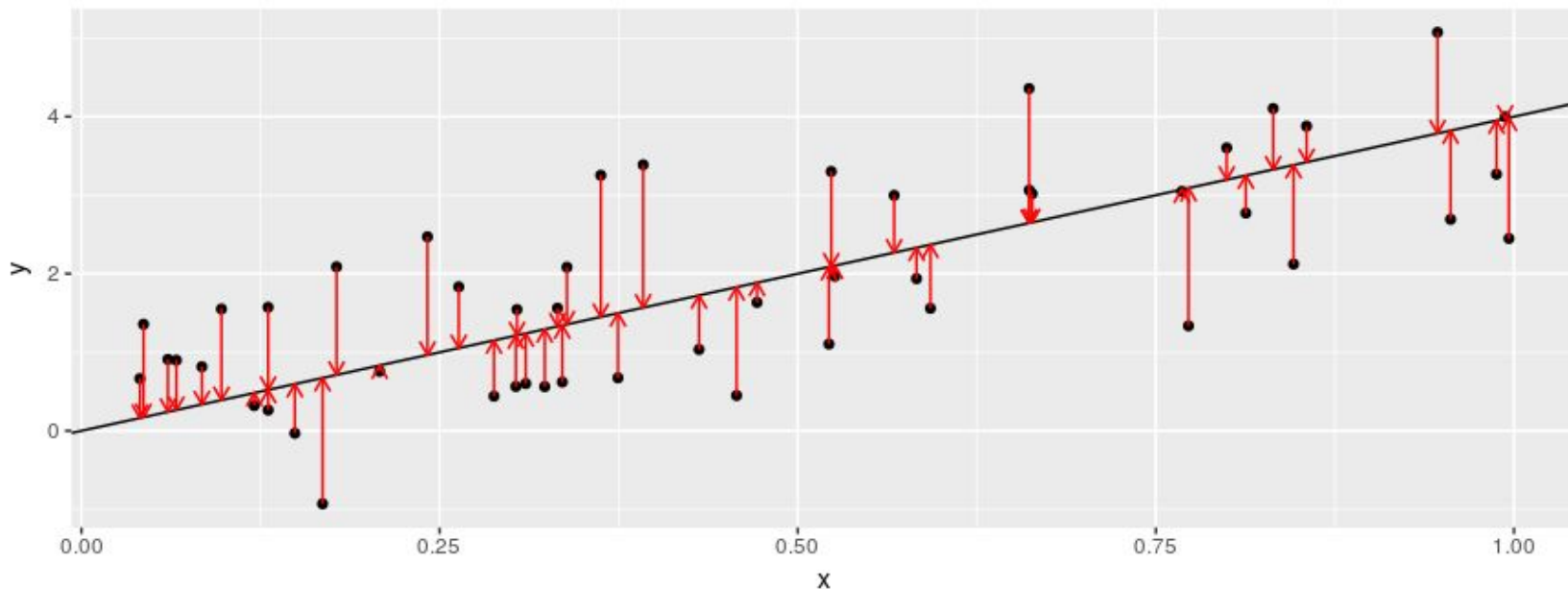
Paul Gustafson

# Regression example: Bayesian leverage scores

Consider the “fitted values”, or “predictions”:

$$\hat{Y} := X\beta$$

As a function of the unknown parameter, these have a posterior distribution.



# Regression example: Bayesian leverage scores

Recall that:

$$\mathbb{E}[\beta|Y] = (\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1} \sigma^{-2}X^T Y$$

Plugging in:

$$\begin{aligned}\mathbb{E}[\hat{Y}|Y] &= \mathbb{E}[X\beta|Y] \\ &= X\mathbb{E}[\beta|Y] \\ &= X(\sigma^{-2}X^T X + \sigma_0^{-2}I_p)^{-1} \sigma^{-2}X^T Y\end{aligned}$$

# Regression example: Bayesian leverage scores

The prediction depends on  $Y$ , and we can calculate the local sensitivity to  $Y$ .

$$\mathbb{E} [\hat{Y}|Y] = X (\sigma^{-2} X^T X + \sigma_0^{-2} I_p)^{-1} \sigma^{-2} X^T Y$$

$$\frac{d}{dY} \mathbb{E} [\hat{Y}|Y] = X \left( X^T X + \left( \frac{\sigma}{\sigma_0} \right)^2 I_p \right)^{-1} X^T$$

$$\Rightarrow \frac{d}{dy_i} \mathbb{E} [\hat{y}_i|y_i] = \left( X \left( X^T X + \left( \frac{\sigma}{\sigma_0} \right)^2 I_p \right)^{-1} X^T \right)_{ii}$$

**This is the (Bayesian) leverage score!**

# Sensitivity is covariance is sensitivity

This is a general result, so let's state it in some generality. Choose a family of distributions indexed by  $t$ :

$$p(\theta|t) \text{ for some } t \in \mathbb{R}^k$$

We're interested in the sensitivity to  $t$  of this posterior expectation:

$$\mathbb{E}_{p(\theta|t)} [\theta] = \int \theta p(\theta|t) d\theta$$

It turns out that:

$$\frac{d\mathbb{E}_{p(\theta|t)} [\theta]}{dt} = \text{Cov}_{p(\theta|t)} \left( \theta, \frac{d}{dt} \log p(\theta|t) \right)$$

# Sensitivity is covariance is sensitivity

$$\frac{d\mathbb{E}_{p(\theta|t)}[\theta]}{dt} = \frac{d}{dt} \int \theta p(\theta|t) d\theta$$

$$= \int \theta \frac{d}{dt} p(\theta|t) d\theta$$

$$= \int \theta \frac{d}{dt} \exp(\log p(\theta|t)) d\theta$$

$$= \int \theta \frac{d}{dt} (\log p(\theta|t)) p(\theta|t) d\theta$$

$$= \int \theta \left( \frac{d}{dt} \log p(\theta|t) - \mathbb{E}_{p(\theta|t)} \left[ \frac{d}{dt} \log p(\theta|t) \right] \right) p(\theta|t) d\theta$$



Score is  
mean zero

$$= \int (\theta - \mathbb{E}_{p(\theta|t)}[\theta]) \left( \frac{d}{dt} \log p(\theta|t) - \mathbb{E}_{p(\theta|t)} \left[ \frac{d}{dt} \log p(\theta|t) \right] \right) p(\theta|t) d\theta$$

$$= \text{Cov}_{p(\theta|t)} \left( \theta, \frac{d}{dt} \log p(\theta|t) \right)$$

(Some regularity conditions are required to differentiate the log probability and to exchange integration and differentiation)



# Idea: MCMC for approximate leverage scores

Recall the Bayesian leverage score:

$$\frac{d}{dy_i} \mathbb{E} [\hat{y}_i | y_i] = i^{th} \text{ leverage score} = \text{Cov}_{p_{\alpha}^x} \left( \hat{y}_i, \frac{\partial \log p(Y, \beta)}{\partial y_i} \right) = \text{Cov}_{p_{\alpha}^x} (x_i \beta, \sigma^{-2} x_i \beta)$$

A step of Hamiltonian Monte Carlo (HMC) requires only the evaluation of:

$$\frac{\partial \log p(Y, \beta)}{\partial \beta} = - (\sigma^{-2} X^T X + \sigma_0^{-2} I_p) \beta + \sigma^{-2} Y^T X + C$$

Each requires only  $O(NP)$  time.

# Idea: MCMC for approximate leverage scores

- Given a matrix  $X$ , generate data according to  $p(Y, \beta)$
- Run HMC to generate  $K$  draws from the posterior  $p(Y|\beta)$
- Calculate the approximate leverage scores using the sample covariance:

$$\bar{\beta} := \frac{1}{K} \sum_{k=1}^K \beta_k$$

$$\frac{d}{dy_i} \mathbb{E} [\hat{y}_i | y_i] = \text{Cov}_{p_{\alpha}^x} (x_i \beta, \sigma^{-2} x_i \beta) \approx \frac{1}{K} \sum_{k=1}^K \sigma^{-2} (\beta_k - \bar{\beta}) x_i x_i^T (\beta_k - \bar{\beta})$$

**The whole procedure should be  $O(NPK)$ .**

# What is variational Bayes?

---



How do you calculate robustness to function-valued perturbations?

---

